

To date:  $Y = f(X) + u$

Find  $\hat{f}(X)$

Simple prediction. ~~Some~~

Sometimes useful in econometrics.

Now! Consider workhorse models for causal inference in econometrics

(1) IV when instruments are known and valid but there are too many.

If we use them all have too many instruments problems

(2) Selection on observables model.

Consistent estimates once include controls, but there are too many controls.

Solution: Select a subset of variables but

(1) Avoid overfitting / overtraining

(2) Include enough but not too many.

Survey article Belloni et al JEP Spring 2014 pp 29-50

# 240F Instrumental Variables & Lasso

(2)

## - Key articles

Full details: Belloni et al. Econometrica 2012 pp 2369-2429

Survey: Belloni et al JEP Spring 2014 pp 39-41

## - Consider

$$\text{model. } y_i = d_i' \alpha + \varepsilon_i$$

↑ At least one component of  $d_i$  is endogenous

Instruments:  $x_i$  satisfy  $E(\varepsilon_i | x_i) = 0$

The optimal instrument with iid errors is

$$D(x_i) = E(d_i | x_i)$$

and leads to IV estimator with variance  ~~$\frac{1}{n} E[D(x_i) D(x_i)']$~~

matrix  $\sigma^2 \left\{ \frac{1}{n} \sum_{i=1}^n E[D(x_i) D(x_i)'] \right\}^{-1}$ .

- Problem: Clearly best to use instruments that approximate  $D(x_i)$   
eg. regress  $d_i$  on series expansion of  $x_i$   
But this may lead to too many instruments so asymptotic theory poor.

~~Solution:~~

Solution:

~~Alternatively~~ pick just a few functions of  $x_i$

as instruments. But how? Lasso

Use Lasso to pick instruments. Then IV as usual.

- $n =$  sample size
- $p =$  # instruments
- $p$  large relative to  $n$  (even  $p > n$ ) due to
  - (1) many potential  $z_i$  as instruments (many instrument)
  - (2) a few  $z_i$  but consider flexible model with powers and interactions. (many series instruments)
- ~~Use the Lasso. For each endogenous regressor~~  
~~Regress the Lasso on the  $z_i$ 's and then~~  
~~to select the subset of~~  
~~the  $p$~~   
 - Let  $\uparrow$  instruments be
 
$$f_i = f(i_1, \dots, i_p)' = (f_1(z_i), \dots, f_p(z_i))'$$
- Approximate sparsity  
 Assume ~~approximate~~ that ~~only~~ ~~only~~  $s$  of these instruments are needed to provide a good approximation to  $D_\ell(z_i) = E[d_{\ell i} | z_i]$  for the  $\ell^{\text{th}}$  endog. variable  
 i.e.  $D_\ell(z_i) = f_i' \beta_{\ell 0} + a_\ell(z_i)$  for  $\ell = 1, \dots, k$ 

$$\max_{1 \leq \ell \leq k} \sum_{j=1}^p \mathbb{1}[\beta_{\ell 0 j} \neq 0] \leq s = o(n)$$
 where approximation error  $a_\ell(z_i)$  is at most  $O_p(\sqrt{s/n})$

LASSO

$$- d_{i,l} = D_i(x_i) + v_{i,l} \quad E[v_{i,l}(x_i)] = 0 \quad l = 1, \dots, k_l$$

$$D_{i,l} = D_l(x_i) = f_l' \hat{\beta}_l \quad l = 1, \dots, k_l$$

$$- \text{Lasso } \hat{\beta}_{l,l} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n (d_{i,l} - f_l' \beta)^2 + \frac{\lambda}{n} \|\vec{\gamma}_l(\beta)\|,$$

where  $\lambda$  is the penalty level

$$\vec{\gamma}_l = \text{diag}(\hat{\gamma}_{l1}, \dots, \hat{\gamma}_{lp}) \quad \text{penalty loadings}$$

$$\text{Ideally } \hat{\gamma}_{lj} = \sqrt{\frac{1}{n} \sum_{i=1}^n f_{ij}^2 v_{i,l}^2} \quad j = 1, \dots, p$$

- Since  $v_{i,l}$  is not observed start with conservative penalty loadings and then plug in the resulting estimator residuals, iterate.

$$\lambda = c 2\sqrt{n} \Phi^{-1}\left(1 - \frac{\delta}{2k_l p}\right)$$

$$\delta = \frac{0.1}{\ln(\max(p, n))}$$

$$c = 1.1$$

- Once lasso picks the  $\hat{\beta}_l \neq 0$   
post lasso does OLS with just these regressors.  
(OLS avoids shrinkage bias, lasso avoids overfitting).

IV Form  $\hat{D}_l(x_i) = f_l' \hat{\beta}_l \quad l=1, \dots, k$

↑  
 Called be either Lasso  $\hat{\beta}_l$   
 or post Lasso  $\hat{\beta}_l$

And  $\hat{D}_i = (\hat{D}_i(x_i), \dots, \hat{D}_{k_e}(x_i), w_i')'$

(Aside what if  $\hat{D}_l(x_i)$  involves  $w_i$  ? )  
 ↑ exogenous regressors

IV  $\hat{\alpha} = \left( \frac{1}{n} \sum_{i=1}^n \hat{D}_i \hat{D}_i' \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \hat{D}_i y_i \right)$

Can use usual asymptotic theory.

With heteroscedastic errors use vce (robust)

$$\hat{V}(\hat{\alpha}) = \left( \frac{1}{n} \sum_{i=1}^n \hat{D}_i \hat{D}_i' \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \hat{D}_i \hat{D}_i' \right)^{-1} \left( \sum_{i=1}^n \hat{D}_i \hat{D}_i' \right)$$

- ① model selection mistakes are small enough to not matter
- ② lasso avoids overfitting.

Application

$$y_{ct} = \alpha_c + \alpha_t + \delta_{ct} + \beta \text{ Takings law} + w_{ct}'\delta + \epsilon_{ct}$$

$\uparrow$  judicial circuit c  
 $\uparrow$  time t  
 Instruments arise due to random assignment of judges

n = 183

p = 147

# of takings appellate decisions that rule that a taking was unlawful

I think they ~~factor~~ <sup>partial</sup> out  $\alpha_c, \alpha_t, \delta_{ct}, w_{ct}'\delta$

So just  $\hat{y}_{ct} = \alpha + \beta \text{ Takings law} + \text{error}$ .

Focus on  $y_{ct}$  = Case-Shiller home price index within circuit c at time t

Eta paper finds		S = 2 instruments	
JEP		S = 1	
$\hat{\beta}$	OLS	2SLS	lasso
$se(\hat{\beta})$	.0152	.0604	.0631
	.0132	.0296	.0249

lasso shooting y z\*

$\uparrow$   
 Chou & Yeh (2010)  
 First stage R<sup>2</sup> = 89.6  
 First stage F = 67.7

240F

## Control Function & Lasso

⑦

Full details Belloni et al. REST <sup>nd</sup> 2014 608-650

Survey Belloni et al JEP Spring 2014 41-45

- Consider partially linear model

$$y_i = d_i \alpha_0 + g(z_i) + \varepsilon_i \quad E(\varepsilon_i | z_i, d_i) = 0$$

If  $z_i$  was low dimension can eg. use Robinson's 1998 differencing estimator.

- Instead there are many  $z_i$ :

In fact  $p$  potential regressors  $x_i = p(z_i)$

- Assume enough sparsity so that only  $s < n$  of  $x_i$  matter,

~~are~~ approximating

$g(z_i)$  by linear combination of them

and  $E(d_i | z_i) = m(z_i)$ .

$$y_i = d_i \alpha_0 + g(z_i) + \varepsilon_i \quad E(\varepsilon_i | z_i, d_i) = 0$$

$$d_i = m(z_i) + v_i \quad E(v_i | z_i) = 0$$

Approximate by linear combinations of  $d_i = p(z_i)$

$$y_i = d_i \alpha + \underbrace{x_i' \beta_g}_{g(z_i)} + r_{gi} + \varepsilon_i$$

$$d_i = \underbrace{x_i' \beta_m}_{m(z_i)} + v_i$$

$x_i = p(z_i)$  has dimension  $p = p_n$  with  $\log p = o(n^{1/3})$

B-splines, dummies, polynomials, interactions, ...

Sparsity

$$\textcircled{1} \quad \|\beta_m\|_0 \leq s \text{ and } \|\beta_g\|_0 \leq s \text{ for } s \ll n$$

↑ # non-zero  $\beta$

$\textcircled{2}$  Approximation errors are small relative to estimation error

$$\sqrt{\sum_i E(r_{gi}^2)} \leq \frac{c}{n} \times \sqrt{s/n} \text{ for some constant } c > 0$$

$$\sqrt{\sum_i E(v_{mi}^2)} = c \times \sqrt{s/n}$$



Double Selection

RF ①  $y_i = x_i' \bar{\beta}_0 + \bar{\tau}_i + \bar{\xi}_i$

$$\bar{\beta}_0 = \alpha_0 \beta_m + \beta_{g_0}$$

②  $d_i = x_i' \beta_m + \tau_{mi} + v_i$

$$\bar{\tau}_i = \alpha \tau_{mi} + \tau_{gi}$$

$$\bar{\xi}_i = \alpha \psi_i + \xi_i$$

sub out  $d_i$  in first eqn.

Apply Lasso to each equation separately!

Similar  $\lambda$  penalty to FITE paper.

Then choose  $x_i$  to be union of  $x_i$ 's from ① & ② (Double selection).

Then do OLS of  $y_i$  on  $d_i$  and these  $x_i$ 's.

Can apply regular asymptotics

Why double selection?

Will sometimes drop a <sup>control</sup> variable in the original equation when  $\beta_g \approx 0$ . But it still could be really important if  $\beta_m \neq 0$  and  $\alpha \neq 0$ .

240F

# Control Function & Lasso

(10)

- Paper section 5 discusses heterogeneous effects.

Also see Max Farrell 2015 J. Econometrics 189 1-23

- Application

Donohue & Levitt Abortion and Crime

$$y_{it} = \alpha a_{it} + w_{it}'\beta + \delta_i + \gamma_{it} + \varepsilon_{it}$$

Crime  
 violent  
 property  
 murder

↑  
 Abortion rate

284 controls plus  
 time effects  
 that are partialled out

n = 48  
T = 12  
nT = 576

or  
n = 50  
T = 12  
nT = 600

Doubly Robust Estimation of ATE

- Max Farrell 2015 J. Econometrics 189 1-23

- Multivalued treatment  $D \in \{0, 1, \dots, J\}$

Generalized treatment score  $p_t(x) = \Pr(D=t | X=x)$

Conditional outcome mean,  $E(Y | D=t, X=x) =$   
function  $\mu_t(x) =$

iid sample  $\{y_i, d_i, x_i\} \quad i=1, \dots, n$

Assume selection on observables

$$\hat{\mu}_t = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{\mathbb{1}(d_i=t)(y_i - \hat{\mu}_t(x_i))}{\hat{p}_t(x_i)} + \hat{\mu}_t(x_i) \right\}$$

1. Fit  $p_t$  &  $\mu_t$  using group lasso applied to multinomial logit & LS
2. Refit  $p_t, \mu_t$  using unpenalized models possibly augmented with controls from theory / prior work

Apply to Dehejia Wahba (1999) like Robins, Rotnitzky (1995)

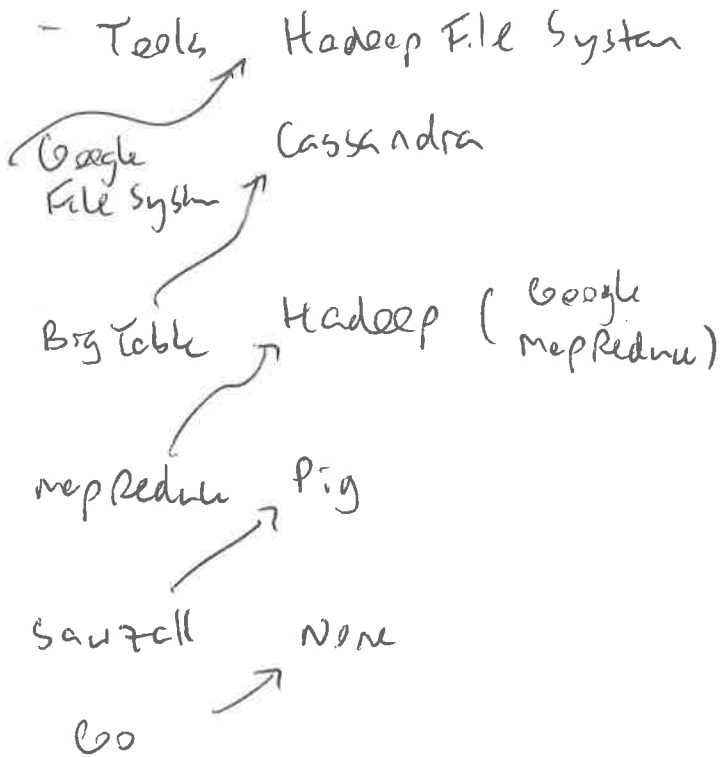
Grouped lasso ESL p. 90 with  $L$  groups ( $y$  is  $L \times 1$ )

$$\min_{\beta} \left( \left\| \begin{matrix} y \\ \sim \end{matrix} - \beta_0 e - \sum_{\ell=1}^L X_{\ell} \beta_{\ell} \right\|_2^2 + \lambda \sum_{\ell=1}^L \sqrt{p_{\ell}} \|\beta_{\ell}\|_2 \right)$$

$\|\cdot\|_2$  denotes Euclidean norm without squaring

$$\|a\|_2 = \sqrt{a'a}$$

- Hal Varian Big Data: New Tricks for Econometrics  
JEP Spring 2014 3-28



Large Files

Table of data that can stretch over many computers

System to access & manipulate data in large data structures

Language to create mapreduce jobs

Language for parallel data processing

SQL queries

Practical  
Big Query

HW  
Drill

- Predicted regression
- Training tests validation
- K-fold CV
- Classification with regression trees

Further work  
Imbens & Athey